

GP-ALPS: Automatic Latent Process Selection for Multi-Output Gaussian Process Models

Pavel Berkovich^{*†}
 Eric Perim^{*}
 Wessel Bruinsma^{*‡}

P.BERKOVICH@CS.UCL.AC.UK
 ERIC.PERIM@INVENIALABS.CO.UK
 WPB23@CAM.AC.UK

1. Introduction

A principled approach to prediction tasks is to choose a statistical model that explains the data. The choice of the *model class* is crucial and has to observe the *bias–variance trade-off*, which motivates the need for principled approaches to selecting the best model class from a set of options. Whilst model selection can be done manually by trial and error, the process tends to consume considerable time and resources and be prone to human biases. Bayesian model selection (MacKay, 1992; Kuo and Mallick, 1998; Rasmussen and Ghahramani, 2001), treats the model class as a *random variable* and computes its posterior distribution. It offers a built-in complexity regulariser, commonly known as Bayesian Occam’s razor, which penalises models whose complexity is excessive or too modest. As a result, Bayesian model selection assigns high posterior probability to model classes whose complexity is “just right”.

Gaussian processes (GPs) are a popular and widely used approach to single-output non-linear regression (Williams and Rasmussen, 2006). They constitute a probabilistic modelling framework that is tractable, modular, and interpretable. GPs can be extended to multiple output and have in this setting successfully been applied to problems as diverse as analysis of neuron activation patterns (Yu et al., 2009), image upscaling (Akhtar et al., 2016), and solar panels’ output prediction (Dahl and Bonilla, 2018). One of the simplest and most widely adopted approach to extend GPs to multiple outputs is to model each output as a linear combination of a collection of shared, unobserved latent Gaussian processes (Wackernagel et al., 1997), henceforth referred to as the Linear Mixing Model (LMM). A pressing issue with this approach is choosing the complexity of the latent space, which constitutes choosing the number of latent processes and their kernels. These choices are typically done manually (Teh and Seeger, 2005; Osborne et al., 2008; Yu et al., 2009), which can be time consuming and prone to overfitting.

In this work, we apply Bayesian model selection to the calibration of the complexity of the latent space. We propose an extension of the LMM that automatically chooses the latent processes by turning off those that do not meaningfully contribute to explaining the data. We call the technique Gaussian Process Automatic Latent Process Selection (GP-ALPS). The extra functionality of GP-ALPS comes at the cost of exact inference, so we devise a variational inference (VI) scheme and demonstrate its suitability in a set of preliminary experiments. We also assess the quality of the variational posterior by comparing

^{*} Invenia Labs, Cambridge, UK

[†] University College London

[‡] University of Cambridge

our approximate results with those obtained via a Markov Chain Monte Carlo (MCMC) approach.

2. Automatic Latent Process Selection (ALPS) for MOGPs

We adopt the following formulation of the Linear Mixing Model:

$$x_j \sim \mathcal{GP}(0, k_j(t, t')), \quad f(t) = Hx(t), \quad y_i(t) \sim \mathcal{N}(f_i(t), \sigma_i^2). \quad (1)$$

In the LMM, $f_i(t) = \sum_{j=1}^m H_{ij}x_j(t)$ is a linear combination of unobserved processes $(x_j)_{j=1}^m$, where we call H the *mixing matrix* and x the *latent processes*. Our approach, named Gaussian Process Automatic Latent Process Selection (GP-ALPS), aims to automatically select those latent processes x_j that meaningfully contribute to the observed signal. It does so by multiplying every x_j by a Bernoulli random variable b_j , which gives the model the ability to exclude x_j from contributing to f : $f(t) = H(x(t) \circ b)$, where \circ denotes the Hadamard product. This approach can be interpreted as a form of drop-out regularisation (Nalisnick et al., 2019) on the latent processes. GP-ALPS also includes a prior over H . In summary, GP-ALPS is given by the following generative model:

$$x_j \sim \mathcal{GP}(0, k_j(t, t')), \quad b_j \sim \text{Bern}(\theta_j), \quad H_{ij} \sim \mathcal{N}(0, s_{ij}), \quad (2)$$

$$f(t) = H(x(t) \circ b), \quad y_i(t) \sim \mathcal{N}(f_i(t), \sigma_i^2).$$

Each of the 2^m possibilities for the vector b identifies a model class, so the prior effectively describes an ensemble of 2^m different models, corresponding to all possible combinations of the latent functions. Another interpretation of GP-ALPS is that the latent processes are various features on which the observed signal can depend, which makes GP-ALPS a method to perform automatic feature selection.

3. Variational Inference Scheme

Let $Y \in \mathbb{R}^{p \times n}$ denote observed data at input locations $t \in \mathbb{R}^n$. Augment the model with *inducing variables* $X^z \in \mathbb{R}^{m \times \ell}$ at input locations $t^z \in \mathbb{R}^\ell$, which are assumed to be sufficient statistics for the latent processes (Titsias, 2009a; Hensman et al., 2013; Nguyen and Bonilla, 2014). To perform inference, we introduce a structured mean-field approximate posterior distribution

$$q(X, X^z, H, b) = p(X|X^z)q(X^z)q(H)q(b)$$

where we choose $q(X^z)$, $q(H)$, and $q(b)$ by minimising the Kullback–Leibler divergence with respect to the true posterior, using stochastic gradient-based optimisation:

$$(q^*(X^z), q^*(H), q^*(b)) = \operatorname{argmin}_{(q(X^z), q(H), q(b))} D_{\text{KL}} [q(X, X^z, H, b) \| p(X, X^z, H, b|Y)]$$

We let $q(X^z)$ be a Gaussian that factorises over the latent processes and $q(H)$ a fully factorised Gaussian. The approximate posterior $q(b)$, however, is troublesome, because b is discrete, which means that we cannot just use the reparametrisation trick (Kingma and Welling, 2013; Titsias and Lázaro-Gredilla, 2015; Rezende et al., 2014). We therefore let $q(b)$ be a continuous relaxation of the Bernoulli distribution called the *concrete distribution*

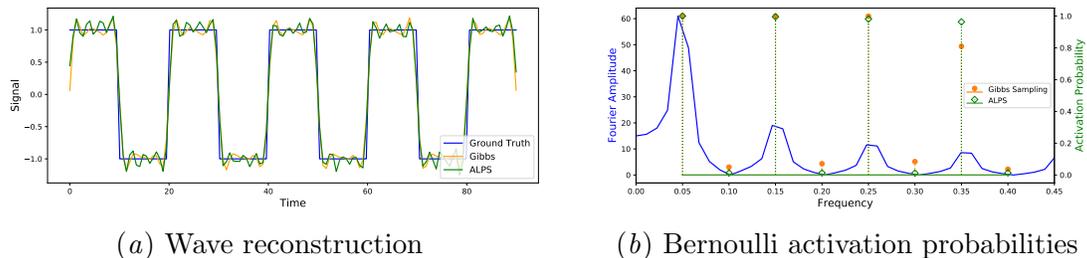


Figure 1: Comparison between GP-ALPS and MCMC for the experiment in Section 4.1.

(Maddison et al., 2016). To compute the ELBO, we also approximate $p(b)$ with the concrete distribution, as the cross-entropy $\mathbb{E}_{q(b)} \log p(b)$ would not be well-defined in the case in which $q(b)$ is continuous and $p(b)$ is discrete (Maddison et al., 2016). For the temperature of the concrete distributions, we use a particular annealing scheme. See Appendix B for a more detailed description of the variational inference scheme.

To assess the quality of the variational approximate posterior, we compare it against Gibbs sampling, which has theoretical guarantees to converge to the true posterior in the infinite time limit. The key insight is that f is bilinear in H , x , which means that $p(X|Y, H, b)$ and $p(H|X, Y, b)$ are tractable (just Bayesian linear regression); and $p(b_j|X, Y, H, b_{-j})$ is tractable because b is discrete. See Appendix C for a more detailed description of the Gibbs sampler.

4. Experimental Results

4.1. Square Wave Decomposition

We first test the model’s ability to select relevant latent processes using a simple example from signal processing. We generate a single ($p = 1$) square wave of frequency $f_{\text{sq}} = 0.05$ Hz, and aim to model it as a linear combination of $m = 8$ latent GPs with linear-periodic kernels with fixed frequencies $f_i = i f_{\text{sq}}$ for $i \in 1, \dots, m$. As can be seen in Figure 1(b), GP-ALPS assigns high Bernoulli activation probabilities to the latents whose frequencies are odd multiples of f_{sq} , which correspond to the peaks in the square wave’s power spectrum. The signal is thus reconstructed quite accurately (Figure 1(a)) using the first 4 terms of the Fourier series. Furthermore, both the activation probabilities and signal reconstruction found by GP-ALPS are quite close to those obtained by sampling the exact posterior via Gibbs sampling, which indicates that the variational posterior approximates the exact one closely, despite the simplifying assumptions that have been made to enable variational inference. Interestingly, the above is true with as few as $\ell = 10$ learnable inducing points.

4.2. Noisy Mixture of Periodic Signals

In this experiment, we test the technique’s ability to perform model selection in the presence of noise, as well as choose between equally good solutions. To generate the data, we start with $m^* = 3$ signals with periods 7, 17 and 23 (Figure 2(a)), then corrupt them with additive Gaussian noise and combine linearly with a fixed matrix $H^* = [I_3 \ Z]^T$ (where

$Z \in \mathbb{R}^{6 \times 6}$), to obtain $p = 9$ outputs (blue in Figure 2(b)). We model the data with GP-ALPS with $m = p = 9$ linear-periodic latents, with periodicities 3, 7, 7, 11, 13, 17, 19, 23 and 23 (note the duplicates), and $\ell = 100$ learnable inducing points.

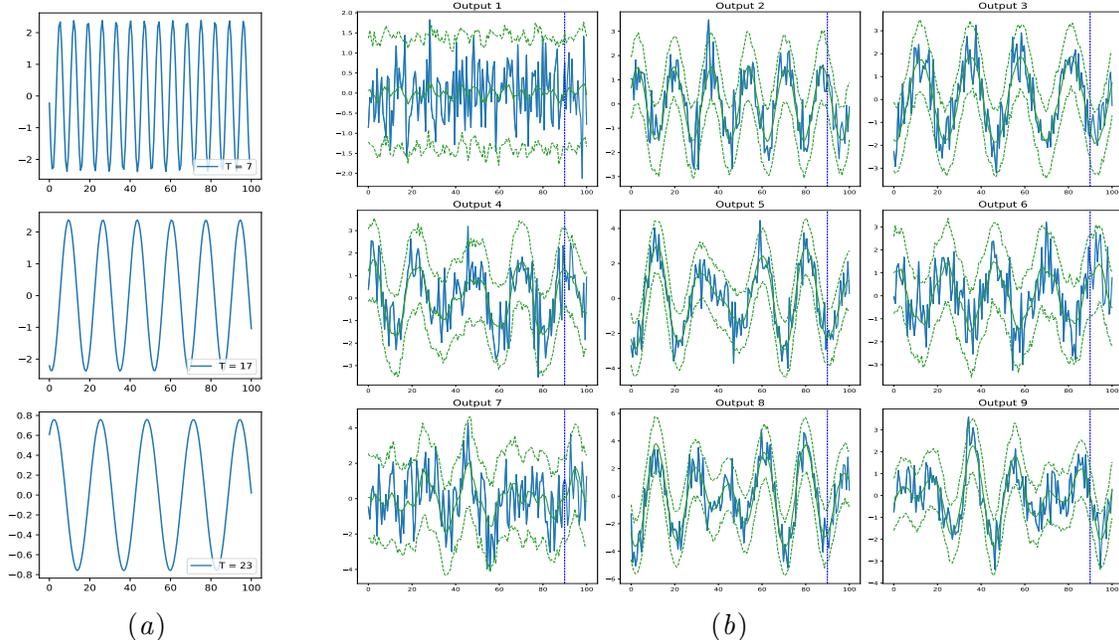


Figure 2: Data generated for the experiment in Section 4.2. (a) Original, noiseless latent signals; (b) outputs (blue) and predictions (green). Vertical line separates training and testing sets.

The predictive densities are shown in green in Figure 2(b), and the trained variational posteriors $q(H)$ and $q(b)$ are shown in Figures 3(a), 3(b), and 3(c). GP-ALPS successfully identifies the frequencies that generated the data. Despite the noise making it virtually impossible to identify the periodicity $T = 7$ visually in the data (Output 1 in Figure 2(b)), the model manages to identify its presence with high degree of certainty. Furthermore, the solution found by GP-ALPS is parsimonious—only one latent is activated for each $T = 7$ and $T = 23$. While, intuitively, one may expect both latents with $T = 7$ (or with $T = 23$) to split the activation probability of those frequencies, this is a “more complex” explanation (either can be on or off) than activating only one (only one can be on or off). By Bayesian Occam’s Razor, we expect that posterior inference tends towards the simpler explanation. Inspecting the element-wise posterior means and variances in H (Figure 2(a)), we note that elements in activated columns are estimated with low-variance Gaussians, as expected, whereas the inactive columns just revert to the standard normal prior.

4.3. Variable Selection in Boston Housing Dataset

Further to the experiments with synthetic data described above, we have employed GP-ALPS to perform variable selection for kernelised ridge regression (KRR), using the Boston

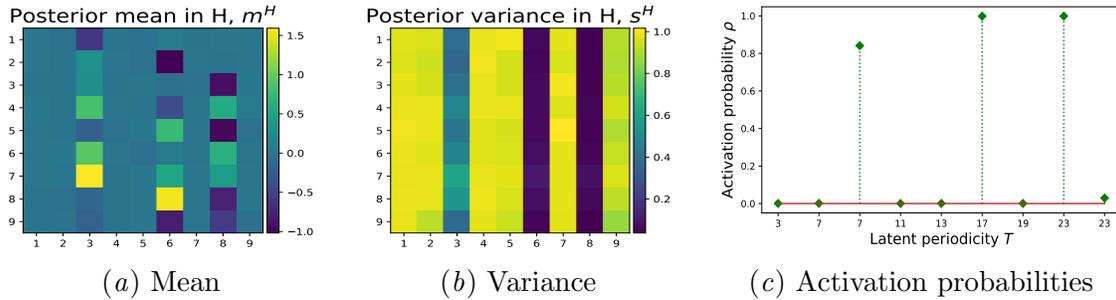


Figure 3: Approximate posteriors from the experiment in Section 4.2.

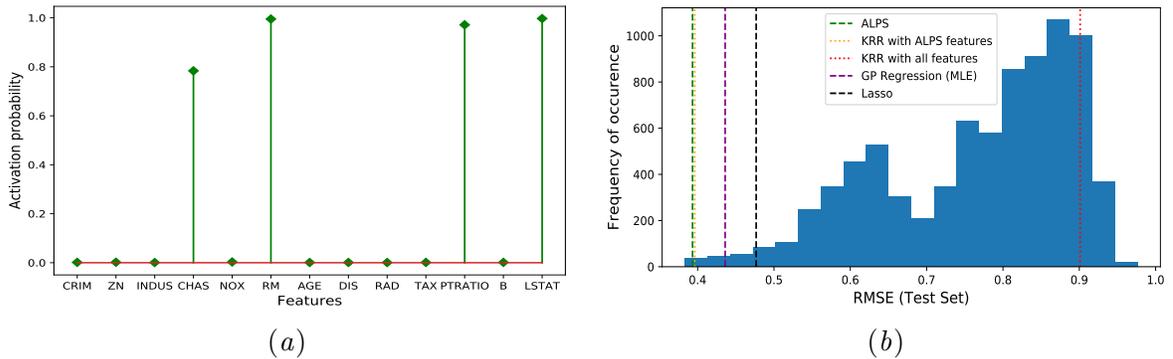


Figure 4: (a) Feature activation probabilities found by GP-ALPS. (b) Comparison between GP-ALPS and all possible 2^{13} kernelised ridge regression (KRR) models on test-set RMSE.

housing dataset¹ as a motivating example. Posterior activation probabilities are shown in Figure 4(a). Comparing the test-set results with all $2^{13} = 8192$ possible linear regression models, we demonstrate that our method performs competitively, ranking within the 0.05% best models, as shown in Figure 4(b). This performance is comparable to the one achieved by KRR using only the features selected by GP-ALPS and superior to the one obtained by carrying regular GP regression with all features, as well as that obtained using Lasso regression. More details can be found in Appendix A.

1. https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_boston.html#sklearn.datasets.load_boston

References

- Naveed Akhtar, Faisal Shafait, and Ajmal Mian. Hierarchical beta process with gaussian process prior for hyperspectral image super resolution. In *European Conference on Computer Vision*, pages 103–120. Springer, 2016.
- Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Springer-Verlag New York, 2006.
- A. Dahl and E. V. Bonilla. Grouped Gaussian processes for solar power prediction. *arXiv preprint arXiv:1806.02543*, 6 2018.
- Amir Dezfouli and Edwin V Bonilla. Scalable inference for gaussian process models with black-box likelihoods. In *Advances in Neural Information Processing Systems*, pages 1414–1422, 2015.
- J. Hensman, N. Fusi, and N. D. Lawrence. Gaussian processes for big data. *arXiv preprint arXiv:1309.6835*, 9 2013.
- D. P. Kingma and M. Welling. Auto-encoding variational Bayes. *arXiv preprint arXiv:1312.6114*, 12 2013.
- Lynn Kuo and Bani Mallick. Variable selection for regression models. *Sankhyā: The Indian Journal of Statistics, Series B*, pages 65–81, 1998.
- David JC MacKay. Bayesian interpolation. *Neural computation*, 4(3):415–447, 1992.
- Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016.
- Eric Nalisnick, José Miguel Hernández-Lobato, and Padhraic Smyth. Dropout as a structured shrinkage prior. In *International Conference on Machine Learning*, pages 4712–4722, 2019.
- Trung V. Nguyen and Edwin V. Bonilla. Collaborative multi-output Gaussian processes. *Conference on Uncertainty in Artificial Intelligence*, 30, 2014.
- M. A. Osborne, S. J. Roberts, A. Rogers, S. D. Ramchurn, and N. R. Jennings. Towards real-time information processing of sensor network data using computationally efficient multi-output Gaussian processes. In *Proceedings of the 7th International Conference on Information Processing in Sensor Networks*, IPSN '08, pages 109–120. IEEE Computer Society, 2008. doi: 10.1109/IPSIN.2008.25.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- Carl Edward Rasmussen and Zoubin Ghahramani. Occam’s razor. In *Advances in neural information processing systems*, pages 294–300, 2001.

- Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. *arXiv preprint arXiv:1401.4082*, 2014.
- Yee Whye Teh and Matthias Seeger. Semiparametric latent factor models. *International Workshop on Artificial Intelligence and Statistics*, 10, 2005.
- Michalis Titsias. Variational learning of inducing variables in sparse gaussian processes. In *Artificial Intelligence and Statistics*, pages 567–574, 2009a.
- Michalis Titsias and Miguel Lázaro-Gredilla. Local expectation gradients for black box variational inference. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 2638–2646. Curran Associates, Inc., 2015.
- Michalis K. Titsias. Variational learning of inducing variables in sparse Gaussian processes. *Artificial Intelligence and Statistics*, 12:567–574, 2009b.
- Hans Wackernagel, Victor De Oliveira, and Benjamin Kedem. Multivariate geostatistics. *SIAM Review*, 39(2):340–340, 1997.
- Christopher KI Williams and Carl Edward Rasmussen. *Gaussian processes for machine learning*, volume 2. MIT press Cambridge, MA, 2006.
- Byron M. Yu, John P. Cunningham, Gopal Santhanam, Stephen I. Ryu, Krishna V. Shenoy, and Maneesh Sahani. Gaussian-Process factor analysis for low-dimensional single-trial analysis of neural population activity. *Advances in Neural Information Processing Systems*, 21:1881–1888, 2009.

Appendix A. Feature Selection in Boston Housing Dataset

We use GP-ALPS to perform feature selection for kernelised linear regression. Our illustrative example is the Boston housing dataset, as provided in `sklearn` (Pedregosa et al., 2011), which contains information about properties in 506 neighbourhoods in Boston, including median value, average number of rooms, average age and some others. The regression task is to predict the median value of a property based on 13 other neighbourhood features.

We model the data using GP-ALPS, whereby each of the $m = 13$ latent processes corresponds to one of the input variables and the latent kernels are radial basis functions (RBF) with unit lengthscale, which is equivalent to kernelised linear regression. The number of inducing points used is $\ell = 100$, and their locations are learnt. GP-ALPS selects 4 out of 13 latent processes (activation probabilities shown in Figure 4(a)) that correspond to variables CHAS (proximity to Charles river), RM (average number of rooms), PTRATIO (average pupil-teacher ratio in local schools) and LSTAT (proportion of population of lower socioeconomic status). Since the size of this data set is comparatively small, it is possible to compare the predictive performance of the model selected by GP-ALPS with that of all the other $2^{13} = 8192$ possible models. Figure 4(b) shows the resulting histogram of test-set root-mean-squared-errors (RMSE) produced by the 8192 kernelised linear regression models. The variable set found by GP-ALPS corresponds to the top-performing 0.05% of the model space. To provide a basis for comparison, we also perform kernelised ridge regression with all 13 variables, GP regression with the kernel comprising a weighted sum of unit-lengthscale RBFs, as well as Lasso regression.

Appendix B. Variational Inference Scheme

As explained in Section 2, the generative model in GP-ALPS explains the observed signal $y(t) \in \mathbb{R}^p$ as a linear-Gaussian transformation of m latent Gaussian processes $x(t) \in \mathbb{R}^m$, multiplied by a vector of Bernoulli variables b . Mathematically, this can be written down as follows:

$$\begin{aligned} x_j &\sim \mathcal{GP}(0, k_j(t, t')), & b_j &\sim \text{Bern}(\theta_j), & H_{ij} &\sim N(0, s_{ij}), \\ f(t) &= H(x(t) \circ b), & y_i(t) &\sim N(f_i(t), \sigma_i^2), \end{aligned}$$

where \circ refers to Hadamard product. Our goal is to compute the posterior over the latent variables, $p(X, H, b|Y)$, but this density is unfortunately intractable, so we resort to variational inference, aiming to find some distribution $q(X, H, b)$ that closely approximates the exact posterior $p(X, H, b|Y)$.

B.1. Analytical formulation

We start by augmenting the latent processes with *inducing variables* $X^z \in \mathbb{R}^{m \times \ell}$ at inducing locations $t^z \in \mathbb{R}^\ell$ (Titsias, 2009a; Hensman et al., 2013; Nguyen and Bonilla, 2014). This construction provides both a meaningful way of summarising the data as part of the posterior on x and an efficient way of scaling to large datasets. With this addition, the evidence lower bound (ELBO) becomes:

$$\mathcal{L} = \mathbb{E}_q \left[\log \frac{p(Y, X, X^z, H, b)}{q(X, X^z, H, b)} \right],$$

which we optimise numerically as in Kingma and Welling (2013). To make this optimisation tractable, we make three important assumptions. We make a structured mean-field assumption, $q(X, X^z, H, b) = q(X, X^z)q(H)q(b)$. We take $q(X, X^z) = p(X | X^z)q(X^z)$, as in Titsias (2009b). For $q(X^z)$ and $q(H)$, we choose fully factorised Gaussians:

$$q(X^z) = \prod_{j=1}^m \mathcal{N}(X_j^z; m_j^z, S_j^z), \quad q(H) = \prod_{i=1}^p \prod_{j=1}^m \mathcal{N}(h_{ij}; m_{ij}^H, s_{ij}^H),$$

where $m^z \in \mathbb{R}^{m \times \ell}$, $S^z \in \mathbb{R}^{m \times \ell \times \ell}$, $m^H \in \mathbb{R}^{p \times m}$, $S^H \in \mathbb{R}^{p \times m}$. Since the Bernoulli distribution does not have a differentiable reparametrisation, we use a continuous relaxation of the Bernoulli distribution for $q(b)$:

$$q(b) = \prod_{j=1}^m \text{Concrete}(b_j; \rho_j),$$

where Concrete is the *concrete* distribution (Maddison et al., 2016). Here $\rho \in [0, 1]^m$. The ELBO can then be re-written as:

$$\mathcal{L} = \underbrace{\mathbb{E}_q(\log p(Y|X, H, b))}_{\stackrel{\text{def}}{=} \mathcal{L}_{\text{ell}}} - \underbrace{D_{\text{KL}}[q(X^z) \| p(X^z)]}_{\stackrel{\text{def}}{=} \mathcal{L}_{\text{kl}}^z} - \underbrace{D_{\text{KL}}[q(H) \| p(H)]}_{\stackrel{\text{def}}{=} \mathcal{L}_{\text{kl}}^H} - \underbrace{D_{\text{KL}}[q(b) \| p(b)]}_{\stackrel{\text{def}}{=} \mathcal{L}_{\text{kl}}^b}$$

Let us consider each of the terms above in turn.

B.2. Expected log-likelihood (\mathcal{L}_{ell})

Start from the full expression for the expected log-likelihood:

$$\mathcal{L}_{\text{ell}} = \mathbb{E}_{p(X|X^z)q(X^z)q(H)q(b)} [\log p(Y|X, H, b)].$$

Since the conditional likelihood inside the expectation does not depend on X^z , we first marginalise it out, similarly to [Dezfouli and Bonilla \(2015\)](#):

$$q(X) = \int dX^z p(X|X^z)q(X^z) = \prod_{j=1}^m \underbrace{\mathcal{N}(X_j; A_j m_j^z, \tilde{K}_j + A_j S_j A_j^T)}_{q(X_j)}$$

where $A_j = K_{cz}^j (K_{zz}^j)^{-1}$ and $\tilde{K}^j = K_{cc}^j - A_j K_{zc}^j$, such that K_{cc}^j , K_{cz}^j and K_{zc}^j are Gram matrices constructed using latent kernel $k_j(\cdot, \cdot)$ on input vectors t and t^z . Adapting Theorem 1 from [Dezfouli and Bonilla \(2015\)](#) to our parametrisation of $q(X^z)$, we then write the ELL as:

$$\mathcal{L}_{\text{ell}} = \mathbb{E}_{q(H)q(b)} \left[\sum_{t=1}^n \mathbb{E}_{q(x_t)} [\log p(y_t|x_t, H, b)] \right],$$

where $q(x_t) = \mathcal{N}(x_t; d_t, S_t^x)$ such that $(d_t)_j = (A_j m_j^z)_t$ and $(S_t^x)_{jj} = \tilde{K}_{tt}^j + (A_j)_t^T S_j^z (A_j)_t$. The final expression for ELL we use is then:

$$\mathcal{L}_{\text{ell}} = \underbrace{\sum_{t=1}^n \mathbb{E}_{q(H)q(b)q(x_t)} [\log p(y_t|x_t, H, b)]}_{\stackrel{\text{def}}{=} \mathcal{L}_{\text{ell}}^{(t)}}$$

whose gradients we compute using the reparametrisation trick ([Kingma and Welling, 2013](#); [Rezende et al., 2014](#); [Titsias and Lázaro-Gredilla, 2015](#)) on variational posteriors $q(x_t)$, $q(H)$ and $q(b)$.

B.3. KL-divergence in activation variables ($\mathcal{L}_{\text{kl}}^b$)

To compute

$$\mathcal{L}_{\text{kl}}^b = -\mathbb{E}_{q(b)} \left[\log \frac{q(b)}{p(b)} \right],$$

we also approximate $p(b)$ with the concrete distribution. Again, we use the reparametrisation trick to compute gradients.

B.4. KL-divergence in inducing variables ($\mathcal{L}_{\text{kl}}^z$)

Both $q(X^z)$ and $p(X^z)$ are block-diagonal multivariate Gaussians, so the KL-term has a closed analytical form:

$$\begin{aligned} \mathcal{L}_{\text{kl}}^z &= -D_{\text{KL}} [q(X^z) \| p(X^z)] = -\sum_{j=1}^m D_{\text{KL}} [q(X_j^z) \| p(X_j^z)] \\ &= -\sum_{j=1}^m D_{\text{KL}} [\mathcal{N}(X_j^z; m_j^z, S_j^z) \| \mathcal{N}(X_j^z; 0, K^j)] \\ &= -\frac{1}{2} \sum_{j=1}^m \left[\text{Tr}[(K^j)^{-1} S_j^z] + (m_j^z)^T (K^j)^{-1} m_j^z - \ell + \log \frac{|K^j|}{|S_j^z|} \right], \end{aligned}$$

so gradients can be computed analytically or by automatic differentiation.

B.5. KL-divergence in mixing matrix ($\mathcal{L}_{\text{kl}}^H$)

Both $q(H)$ and $p(H)$ are also diagonal multivariate Gaussians, so the KL-term is simply

$$\begin{aligned} \mathcal{L}_{\text{kl}}^H &= -D_{\text{KL}} [q(H) \| p(H)] = -\sum_{i=1}^p \sum_{j=1}^m D_{\text{KL}} [q(H_{ij}) \| p(H_{ij})] \\ &= -\sum_{i=1}^p \sum_{j=1}^m D_{\text{KL}} [\mathcal{N}(H_{ij}; m_{ij}^H, s_{ij}^H) \| \mathcal{N}(H_{ij}; 0, s_{ij})] \\ &= -\sum_{i=1}^p \sum_{j=1}^m \left[\frac{1}{2} \log \frac{s_{ij}}{s_{ij}^H} + \frac{s_{ij}^H + (m_{ij}^H)^2}{2s_{ij}} - \frac{1}{2} \right], \end{aligned}$$

which, again, can be differentiated analytically or using automatic differentiation.

B.6. Summary of the optimisation problem

All in all, the variational objective we aim to maximise is:

$$\mathcal{L} = \mathcal{L}_{\text{ell}} + \mathcal{L}_{\text{kl}}^b + \mathcal{L}_{\text{kl}}^z + \mathcal{L}_{\text{kl}}^H.$$

Observing that \mathcal{L}_{ell} is a sum over data points, and that other terms do not depend on observations Y , the ELBO will also be a sum over data points:

$$\mathcal{L} = \sum_{t=1}^n \mathcal{L}^{(t)} = \sum_{t=1}^n \left[\mathcal{L}_{\text{ell}}^{(t)} + \frac{1}{n} (\mathcal{L}_{\text{kl}}^b + \mathcal{L}_{\text{kl}}^z + \mathcal{L}_{\text{kl}}^H) \right],$$

thus, the objective is amenable to stochastic gradient-based optimisation, which is helpful for scaling the model to large datasets.

B.7. Temperature of the concrete distributions

For the temperature of the concrete distributions in $q(b)$, we use the following annealing scheme:

$$T(n, N) = 0.66 + (10.0 - 0.66) \exp\left(-\frac{(n - 0.75N)^2}{0.083^2 N^2}\right)$$

where n is the current iteration and N is the total number of iterations. This annealing scheme is visualised in Figure 5. The idea behind the temperature starting low is that the rest of the parameters can be optimised before latent processes start being dropped out. As for the temperature parameter in the continuous relaxation of $p(b)$, it is chosen to be $1/2$, as in Maddison et al. (2016).

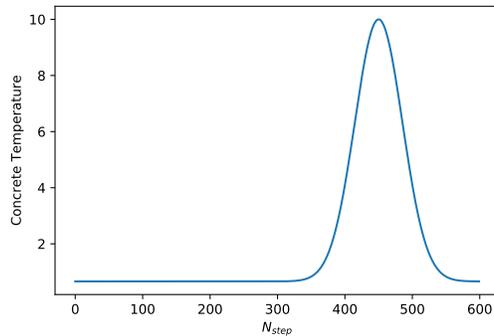


Figure 5: Visualisation of the annealing scheme for the temperature of the concrete distributions

Appendix C. Markov Chain Monte Carlo

This appendix summarises the derivations of the conditionals needed to perform Gibbs sampling from the intractable exact posterior of GP-ALPS. As stated in Section 3, the following three conditionals are of interest:

$$p(X|Y, H, b), \quad p(H|X, Y, b), \quad \text{and} \quad p(b_j|X, Y, H, b_{-j}).$$

Let us consider each of them in turn.

C.1. $p(X|Y, H, b)$

Start by writing down the Bayes' theorem:

$$p(X|Y, H, b) \propto p(Y|X, H, b) p(X) = \mathcal{N}(y; H_b x, S_b) \mathcal{N}(x; 0, K),$$

where $y = \text{vec}(Y)$, $x = \text{vec}(X)$, $H_b = (H \text{diag}(b)) \otimes I_n$, $S_b = \text{diag}(\sigma^2) \otimes I_n$ and K is the block-diagonal multi-output kernel matrix. Note that the above is a Bayesian linear regression problem with a Gaussian prior, so the posterior form is well-known (Bishop, 2006, p. 93):

$$p(x|y, H_b) = \mathcal{N}(x; S_x H_b^T S_b^{-1} y, S_x),$$

where $S_x = [K^{-1} + H_b^T S_b^{-1} H_b]^{-1}$.

C.2. $p(H|X, Y, b)$

As before, writing down the Bayes' theorem:

$$\begin{aligned} p(H|Y, X, b) &\propto p(Y|H, X, b) p(H) = \left[\prod_{i=1}^p p(y_i|X, h_i, b) \right] \left[\prod_{i=1}^p p(h_i) \right] \\ &= \prod_{i=1}^p \left[p(h_i) p(y_i|X, h_i, b) \right] = \prod_{i=1}^p \left[\mathcal{N}(h_i; 0, s_{ij}) \mathcal{N}(y_i; X^T \text{diag}(b) h_i, \sigma_i^2 I_n) \right], \end{aligned}$$

where $h_i \in \mathbb{R}^m$ is i^{th} row of H , and $y_i \in \mathbb{R}^n$ is the i^{th} output. The above amounts to p independent Bayesian linear regression problems, so, as before, the posterior form is well-known (Bishop, 2006, p. 93):

$$p(h_i|y_i, b, X) = \mathcal{N}(h_i; S_h \text{diag}(b) X y_i / \sigma_i^2, S_h),$$

where $S_h = [\frac{1}{s} I_m + \frac{1}{\sigma_i^2} \text{diag}(b) X X^T \text{diag}(b)]^{-1}$.

C.3. $p(b_j|X, Y, H, b_{-j})$

Writing down the Bayes theorem for each b_j :

$$\begin{aligned}
 p(b_j|Y, X, H, b_{-j}) &\propto p(Y|X, H, b) p(b_j) = p(b_j) \left[\prod_{i=1}^p \prod_{t=1}^n p(y_{ti}|b, x_t, h_i) \right] \\
 &\propto \exp \left[b_j \log \frac{\theta_j}{1 - \theta_j} \right] \prod_{i=1}^p \prod_{t=1}^n \exp \left[- \frac{(y_{ti} - \sum_{k=1}^m h_{ik} b_k x_{tk})^2}{2\sigma_i^2} \right] \\
 &= \exp \left[b_j \log \frac{\theta_j}{1 - \theta_j} \right] \prod_{i=1}^p \prod_{t=1}^n \exp \left[- \frac{1}{2\sigma_i^2} (y_{ti}^2 - 2y_{ti} \sum_{k=1}^m h_{ik} b_k x_{tk} + (\sum_{k=1}^m h_{ik} b_k x_{tk})^2) \right] \\
 &= \exp \left[b_j \log \frac{\theta_j}{1 - \theta_j} \right] \prod_{i=1}^p \prod_{t=1}^n \exp \left[- \frac{1}{2\sigma_i^2} (y_{ti}^2 - 2y_{ti} \sum_{k \neq j} h_{ik} b_k x_{tk} - 2y_{ti} h_{ij} b_j x_{tj} \right. \\
 &\quad \left. + (\sum_{k \neq j} h_{ik} b_k x_{tk})^2 + 2h_{ij} b_j x_{tj} + (h_{ij} b_j x_{tj})^2) \right] \\
 &\propto \exp \left[b_j \log \frac{\theta_j}{1 - \theta_j} \right] \exp \left[- \frac{b_j}{2} \sum_{i=1}^p \frac{1}{\sigma_i^2} (2h_{ij} \sum_{t=1}^n \sum_{k \neq j} h_{ik} b_k x_{tj} x_{tk} \right. \\
 &\quad \left. + h_{ij}^2 \sum_{t=1}^n x_{tj}^2 - 2h_{ij} \sum_{t=1}^n x_{tj} y_{ti}) \right] \\
 &= \exp \left[b_j \left(\log \frac{\theta_j}{1 - \theta_j} + c_j \right) \right] = \text{Bern} \left[b_j; \frac{\theta_j e^{c_j}}{e^{c_j} + \theta_j} \right],
 \end{aligned}$$

where

$$c_j = -\frac{1}{2} \sum_{i=1}^p \frac{1}{\sigma_i^2} (2h_{ij} \sum_{t=1}^n \sum_{k \neq j} h_{ik} b_k x_{tj} x_{tk} + h_{ij}^2 \sum_{t=1}^n x_{tj}^2 - 2h_{ij} \sum_{t=1}^n x_{tj} y_{ti}).$$